# SYSTEM AND METHOD FOR MULTIDIMENSIONAL EXTENSION OF DATABASE INFORMATION

# CROSS-REFERENCE TO RELATED APPLICATION

[0001] The subject matter of this application is related to the subject matter of U.S. Provisional Application Serial No. 60/498,283 filed August 28, 2003, from which application this application claims priority.

#### FIELD OF THE INVENTION

[0002] The invention relates to the field of information technology, and more particularly to techniques for generating multidimensional extensions to large-scale medical or other data to permit more efficient searching, data mining and other operations, such as on a clinical or other database.

# **BACKGROUND OF THE INVENTION**

[0003] The advent of powerful servers, large-scale data storage and other information infrastructure has spurred the development of advanced data warehousing and data mining applications. Standard query language (SQL) engines, on-line analytical processing (OLAP) databases and inexpensive large disk arrays have for instance been harnessed in financial, scientific, medical and other fields to capture and analyze vast streams of transactional, experimental and other data. The mining of that data can reveal sales trends, weather patterns, disease epidemiology and other patterns not evident from more limited or smaller-scale analysis.

[0004] In the case of medical data management, the task of receiving, conditioning and analyzing large quantities of clinical information is particularly challenging. The sources of medical data, for instance, may include various independent hospitals, laboratories, research or

other facilities, each of which may generate data records at different times and in widely varying formats. Those various data records may be pre-sorted or pre-processed to include different relationships between different fields of that data, based upon different assumptions or database requirements. When received in a large-scale data warehouse, the aggregation of all such differing data points may be difficult to store in a physically or logically consistent structure. Data records may for instance contain different numbers or types of fields, which may have to be conformed to a standard format for warehousing and searching.

[0005] Even when conditioned and stored, that aggregation of data may prove difficult to analyze or mine for the most clinically relevant or other data, such as those indicating a disease outbreak or adverse reactions to drugs or other treatments. That is in part because the data ultimately stored or accessed for reports may only contain or permit relationships between various parts of the data defined at either the beginning or end of the data management process. That is, the data may reflect only those relationships between different fields or other portions of the data which are defined and embedded by the original data source, or which an end user requests in a query for purposes of generating a report. Relying on source-grouped data is a rigid approach which may omit desired relationships, while relying on back-end queries may tax the OLAP or other query engine being used. Other challenges in receiving, storing and analyzing large-scale medical and other data exist.

#### **SUMMARY OF THE INVENTION**

[0006] The invention overcoming these and other problems in the art relates in one regard to a system and method for multidimensional extension of database information, in which one or more data sources may communicate clinical or other data to network resources including a data enhancement layer before ultimate storage in a data warehouse or other storage facility.

The data enhancement layer along with associated components may prepare and extend the constituent data sets into logical structures reflecting meaningful groupings of the data not present in the raw data source. These multidimensional groupings may likewise be performed before an end user accesses the data warehouse or executes a search. According to embodiments of the invention in one regard, the analytics available to the end user may therefore be more powerful and flexible because they can encompass a greater range of possible groupings and queries. Queries and reports may be made more efficient because potential relationships between data and data attributes may be pre-grouped and stored.

# **BRIEF DESCRIPTION OF THE DRAWINGS**

[0007] The invention will be described with reference to the accompanying drawings, in which like numbers reference like elements.

[0008] Fig. 1 illustrates an overall network architecture in which an embodiment of the invention may operate.

[0009] Fig. 2 illustrates an example source record, of a type which may be processed according to embodiments of the invention.

[0010] Fig. 3 illustrates a logical diagram of a hierarchical grouping, which may be processed according to embodiments of the invention.

[0011] Fig. 4 illustrates an example enhanced multidimensional data grouping, which may be generated according to embodiments of the invention.

[0012] Fig. 5 illustrates a dimensional diagram of data organization, according to embodiments of the invention.

[0013] Fig. 6 illustrates a diagram of the generation of physical storage structures, according to an embodiment of the invention.

[0014] Fig. 7 illustrates a flowchart of overall processing according to an embodiment of the invention.

# **DETAILED DESCRIPTION OF EMBODIMENTS**

[0015] An illustrative environment in which an embodiment of the invention may operate is shown in Fig. 1, in which a data source 102 may communicate clinical medical or other data via a data enhancement layer 110 and other networked components to a transactional data store 130 and ultimately to a searchable set of datamarts 112 for analytic processing. The data source 102 may be or include a medical or other site or facility, such as a hospital, laboratory, university, a military, government or other installation which may capture and store clinical and other data regarding patients, diagnoses, treatments and other aspects or outcomes of medical tests and other medical or other encounters or events.

The data source 102 may transmit one or more source records 118 containing that clinical or other information via a network connection, such as the Internet, local area network (LAN), virtual private network (VPN) or otherwise to a staging database 104, for intermediate storage or processing before being communicated further in the storage service chain. The data source 102 may for instance transmit the source records 118 on a fixed or periodic basis, such as one time per day, week or month, or on a variable or episodic basis, such as when a given amount of data is accumulated, a clinical trial is completed or otherwise.

[0017] Fig. 2 illustrates one example format of one or more of source records 118, in which incident, encounter or other data such as patient identifying information, doctor or other provider identifiers, date fields, diagnostic codes, test results fields and other information may be recorded. In embodiments, source records 118 may also be or include compound records, records which contain links to other records, or other content, formats or functionality.

The staging database 104 may receive the source records 118 and assemble and temporarily or permanently store that data for further transmission and processing. In embodiments staging database 104 may prepare the set of source records 118 for physical storage or logical formatting necessary for downstream warehousing or analytics. As illustrated, according to embodiments of the invention staging database 104 may communicate the source records 118 to a conditioning engine 106 for those purposes. Conditioning engine 106 may be or include, for instance, a server which parses the source records 118 to conform to OLAP or other standards.

[0019] Once any conditioning has been carried out, the source records 118 may be communicated to the data enhancement layer 110 for further processing before committing the source records 118 to permanent or other storage. According to embodiments of the invention in one regard, the data enhancement layer 110 may be or include a server with associated electronic, hard or optical disk storage and other computing, storage or network resources configured or programmed to manipulate source data records, identify or resolve relationships between data components, and store resulting multidimensional groupings to datamarts or elsewhere for analytic processing and other purposes.

[0020] More specifically, according to embodiments of the invention the creation, maintenance and extensions of data relationships that are both hierarchical and multidimensional in nature may be supported and extended via data enhancement layer 110 and other components. According to embodiments of the invention, a set of canonical rules 120 may be used to detect and develop relationships between data or attributes of subject data. The rules 120 may for instance represent or include data pairings which tend to indicate a relationship of interest, such as a causal or correlated relationship. The resulting relationships detected using rules 120, which

may not have been present in or specified by the original data source 102, may then in turn be embedded into or used to build a resulting enhanced data grouping 122, which may be stored to a transactional data store 130 and ultimately made available for searching by end users and others. Among other things, the pre-generation of enhanced data grouping 122 whose cubic or other representation may already include ordered rows, columns, layers or other structures which associate meaningful variables or sets of variables together may enhance to power and efficiency of end user analytics. According to the invention in one regard, the performance of query engines using SQL constructs may for example improve because computationally expensive "join", "group-by" or other operations may be unnecessary.

The dimensions, number of axes, layers or other characteristics of enhanced data grouping 122 may extend beyond the nominal dimensions of the source records 118, aggregations of those records or other raw or original data. The resulting enhanced data grouping 122 may also be specific to or dependent on the original source content, which can be further aggregated into larger identified groups to produce meaningful analytics. According to the invention in one regard, the enhanced data grouping 122 may in embodiments embed or reflect relationships developed between attributes of data, rather than strictly the data values themselves, making manipulation of rules 120 more efficient and storage of enhanced data grouping 122 more economical. It may be noted that the dimensions of enhanced data grouping 122 may in general be unconstrained or freely selected, but may be chosen or changed to conform to particular data models used.

[0022] Due to the open nature of grouping strategies according to the invention in one regard, at least three types of relationships can be detected in the data enhancement layer 110 using rules 120 and other resources. Those types include known, derived and inferred

relationships. The data enhancement layer 110 and other platform components may for one measure known relationships between data elements, such as those embedded in the original data source 102. According to the invention in another regard, ad hoc querying using a manual process may be secondly employed to derive relationships that are not currently recognized or measured, but which may be revealed after interrogating a data store.

[0023] Data mining and analytics according to the invention in another regard can likewise be used to infer a third type of relationship, namely grouped relationships based on statistical quantification, outcomes, measurement and other factors. Following substantiation, particular relationships may pass through a grouping and into the transactional data store 130 or other warehouse environment to populate solution set scenarios supporting analysis based on forecast, hidden or other relationships. Inferred groups may be automatically created based on statistical quantification, allowing an end user to pinpoint correlations between events, such as for example drugs, dosings, procedures, timing of events etc. and outcomes such as extended length of stays, mortality, complications, infections etc. that the end user or facility was not aware of or had not predicted.

[0024] As noted, a conventional approach to data warehousing is to retain the relationships of the data source 102, and if any new relationships are needed, to create those relationships in that source and then extract the relationships into the warehouse facility. If the data groupings necessary for analytics can not be accommodated in the original data source 102, the general conventional approach is to then create them at a back-end or presentation layer through a querying and reporting tool. Due to the complexities of some large-scale data stores, and of health care data in particular, compared for instance to data warehouses in retail, banking

or manufacturing industries, these approaches may not accommodate the analytic demands of end users.

[0025] Addressing these and other disadvantages of a source-driven approach, according to the invention in one regard the data source 102 again may communicate the source records 118 to the data enhancement layer 110 and transactional data store 130 to generate data enhancements including extended or derived groupings not present in the original source records 118. According to the invention in one respect, the data enhancement layer 110 and transactional data store 130 may use the attributes of the original data from source records 118 themselves to define extended dimensions, develop or apply rules 120, grouping configurations and additional element attributes to generate enhanced data grouping 122.

There are at least two types of potential data groupings for extension and other purposes, namely hierarchical and multidimensional. As illustrated in Fig. 3, a hierarchical grouping 132 is a logical structure that uses ordered levels as a means of organizing data. This logical structure is made up of levels, parent and children. A level is a position in a hierarchy, a parent is a value at the level above a given value in a hierarchy and a child is the value at the level under a given value in a hierarchy. This grouping scheme may be used to define a data grouping or aggregation in a hierarchical structure, although it may be noted that in cases a grouping may be generated based on one or more dimension attributes that do not have an ordered relationship. For example, a general ledger account of "Net Income" may be made up of accounts that are allocated to net income. In order to accommodate this variety of data, users may need to be able to specify which sub-accounts "roll-up" into the main account. Thus, a general ledger summary account (e.g. total labor expense) would need to be a hierarchical account that is the parent to a series of children sub-accounts at different levels. This logical

structure may then also be required to specify common costing allocation processes. According to the invention in one regard, data source 102 may be or include source records 118 which are organized in a hierarchical fashion.

As schematically illustrated in Fig. 4, enhanced data grouping 122 may among [0027]other things rely upon a further type of grouping, multidimensional grouping, in a separate logical structure which may be used to create or extend the new dimensions that are based on the values of a series of dimensions or other attributes. Multidimensional groupings may be arranged, for example, as a cube in 3-space in which individual columns, rows and layers reflect different attributes, variables or other quantities or objects. For example, the multidimensional group for the service line of "Cardiology" could be determined as the patients that have the encounter type of inpatient, age greater than 18, diagnosis codes 390.0 – 459.9, physician specialty of cardiologist and a particular nurse unit. This group may then be used to analyze different aggregations for this series of dimensions. According to the invention in another regard, the groupings generated in data enhancement layer 110 may facilitate the analysis of the entire group (represented by the whole cube), one side (A1a - C3a), one column (A1a - A3a), one row (A1a - C1a), one attribute (A1) or other aspects of the enhanced data grouping 122. An attribute can be thought of as an object of reference, either a dimension or fact (modality of Additionally, multidimensional groupings in general and the enhanced data reference). groupings 122 generated according to the invention in particular may have the ability to establish cross-relationships.

[0028] That is, dimensions grouped as members of one group can be grouped as members of another. For example, a physician could be grouped to the both the specialty of "Oncologist" and "Internal Medicine". The driving variable which determines which specialty

the data is grouped to are the values that make up that multidimensional grouping. In other words, the "Oncologist" specialty for this physician may have a different series of values than the "Internal Medicine" specialty. As patient activity occurs for this physician, the combination of the values may then dictate which group may be populated.

By having the ability to accommodate both hierarchical and multidimensional data, embodiments of the invention may support analytics that utilize both types of groups. For example, an end user could use the multidimensional group of "Cardiology" and the hierarchical group of "Net Income" to evaluate the net income that was generated by the cardiology service line, using a single report or analytic tool.

[0030] To accomplish these and other results, according to the invention the data enhancement layer 110 may first acquire and represent the dimensional attributes from one or more data source 102. As illustrated in Fig. 5, once the source records 118 or other original data are acquired, relationships may be defined in or using a schematic physical structure 134 through the application of the following equation:

Equation 1		
		$X^n, Y^n, Z^0 = i > r$
where		
	X, Y =	Dimension attributes
	Z =	Fact reference to transactional activity
	i _> r =	Intersection functionally determines result set.

[0031] Fig. 6 illustrates a diagram of physical structures embedded or used to organize or store source records 118 at various stages of processing according to embodiments of the invention, once those record are acquired. As illustrated, data may be collected from one or more data source 102, such as a health system or company or various other clinical or other facilities or sources. The source records 118 delivered by data source 102 may contain diverse fields or

components, including source facts 126 such as encounters, orders, clinical events and other personal, medical, administrative and other data. The data delivered by data source may likewise include or have associated with it dimensions 128, defining or related to multidimensional cubic or other representations of the data. The source records 118 along with source-defined groupings, source facts 126, dimensions 128 and other data may be assimilated into the process of generating rules 120 by which the aggregate of source data records may be extended by multidimensional groupings, for instance to associate clinically related variables in the same column, row or other space.

[0032] As illustrated in that figure, the generation of rules 120 may be performed or aided by a translation matrix 124, which may be independent from, augment or be part or data enhancement layer 110. In embodiments the rules 120 may be generated by processing the results of prior analytics, by predefined groupings, by automated detection of event or other correlations, or by other techniques. It may be noted that physical structures such as hard disk partitioning of large data arrays may mirror the organization of enhanced data grouping 122 or other components or aspects of data stored to the set of datamarts 112, or other resources.

[0033] After the enhanced data grouping 122 has been defined in a physical storage structure or otherwise, that grouping may be implemented into the transactional data store 130 and ultimately delivered to an appropriate one or more of the set of datamarts 112, for example using further rule translations and thus making the enhanced data grouping 122 transparently available to the end user.

[0034] The following is an example of variables which may be used to generate an enhanced data grouping 122 for the service line of "Cardiology", according to embodiments of the invention:

#### Example 1

X =	DW_Encounters E DW_Encounters_Nomenclature N
Y =	E.Encounter Type – Inpatient E.Age > 18 E.Physician – Cardiologist N.Diagnosis Code 390.0 – 459.9
Z=	# of Admits > 0
i_>r=	Cardiology

[0035] In order to accommodate this flexible approach to groupings which is not source-constrained, the data enhancement layer 110 and other components of the invention may manage and represent fact, dimension and attribute data within at least either a hierarchy or a multidimensional data strategy. In that strategy or implementation both a logical and physical representation of rules and data may be used. Embodiments of the invention may thus represent transactional data elements without the use of core activity data itself, but instead rely solely on the attributes of fact and dimensional data.

[0036] In this regard, the source records 118 and the constituent data may themselves define or be used to define dimensions, rules 120, grouping configurations, element attributes and other criteria used to generate enhanced data grouping 122. As illustrated for example in Fig. 3, the physical structure of the data may identify at least source data relationships, source data element relationships, source data attribute relationships, source data aggregation and source data consolidation, among other things.

[0037] According to the invention in another regard, the management of the mapping of results space to input space may be accomplished by applying "soft data" strategies known to persons skilled in the art. The so-called Soft Data Theorem for instance may be used to take

CRNC.107552

advantage of the fact that dimensions have an inherent hierarchy of determinant data structures and variables, which can be exploited to assist in the generation of enhanced data groupings 122. This approach concentrates on representing data through selective attribute representations of reference data. The technique enables, among other things, the ability to manage multi-dimensional cross-relationships, to relate different levels of aggregation, to relate data at varying granularities, relate data at varying perspectives (cubes, either hyper or multi), density increases at higher consolidation levels and the ability to manage results space to input space for both system-defined and user-defined values.

[0038] The grouping strategies employed by the invention may provide logical aggregation through combining attributes of multiple dimensions 128 that define one group of fact records, as opposed to a physical aggregation that requires schema and foreknowledge of the required dimensions and facts from data source 102 or otherwise. Among other advantages, this may simplify queries by allowing end users to group multiple conditions into a "super group", enhance query performance by reducing the number of SQL or other joins required, allow site-specific dimension groupings, and again enable a common grouping strategy for disparate sources of data.

[0039] As noted, the constructs under which data groupings are applied to generate enhanced data grouping 122 support normalized rules 120 of functional determination. At least three of rules 120 may be fundamental and represent existing data behavior that are determinate, possess normalized relationships and may be inherent to other derivatively-defined relationships. Those three rules among rules 120 are identified as:

TABLE 1

Reflective Rule:	X contains Y, then X -> Y
Augmentation Rule:	${X->Y}$ implies $XZ -> YZ$
Transitive Rule:	${X->Y, Y->Z}$ implies $X->Z$

Other types or classes of inference or other rules may be included within rules 120, including for

example:

# TABLE 2

Resolution: For all clauses C, D and variables A,

 $(C \vee A) (*A \vee D)$ 

(C v D), in which C v D is said to be resolvent, A is a resolved atom.

Factoring: For all clauses C and variables A,

(C v A v A)

(C v A), referred to as C v A factor.

Other types or classes of rules may be used.

[0040] Further, inferential relationships may be hypothesized to have leverage where inferences would exist through possible approximation and differential substantiation. Additional insight into direction of vector(s) intersecting with opposing planes through fact activity supported through the following additional instances of rules within rules 120:

# TABLE 3

Projection Rule:

 $\{X->YZ\}$  implies X->Y

Union Rule:

 $\{X->Y, X->Z\}$  implies X->YZ

Pseudo-Transitive Rule:

 $\{X->Y, WY->Z\}$  implies WX->Z

[0041] These and other rules 120 such as those in the following table represent object and data-related relationships supported within schema, structure and query definitions, such as those supported or required by SQL, OLAP or other data platforms.

# TABLE 4

Base Relations: Physical SQL base tables. Otherwise known as the real relations. These real relations are defined by the physical data warehouse structure.

Views: The virtual relations. A named, derived relation. May also exist as logical layer. Views are defined at the database layer.

Snapshots: A real, not virtual, named derived relation showing the status of an entity at a point in time.

Query Results: The final output relation from a specified query. It may not be named and has no permanent existence. Results can be defined through solution sets or ad hoc query activity.

[0042] According to embodiments of the invention in another regard, additional requirements may arise due to differences that may exist between online transaction processing (OLTP) and OLAP implementations. Transforming OLTP data to an acceptably performing OLAP system may require a number of functionalities.

[0043] Those intermodal functionalities may include:

# TABLE 5

Ability to Merge All Data related to specific items from multiple OLTP systems.

Ability to Resolve Differences in encoding between the different OLTP systems.

Ability to Match Common Data from disparate systems, even data with inconsistencies.

Ability to Convert Different Data types in each OLTP system to a single OLAP type.

Ability to Select Column Data in the OLTP system are not relevant to an OLAP system.

Ability to Absorb Input Data not strictly limited to centrally located OLTP systems.

Ability to Scrub Data - Address inconsistencies to modeled data and process structures. Inconsistencies have to be addressed before data can be loaded into a warehouse for use.

Ability to represent Aggregate Data Relationships notwithstanding transaction details.

Ability to Optimize Aggregate Performances using "Modular Fact Granularities".

Ability to Organize Data in Cubes- Since dimensional attributes are stored in structures designed to represent actual reference data, which already exist in multi-dimensional cube organizations to support analytics, transformation may be achieved through rules structures.

Ability to represent Meta Data Objects in OLTP databases, cubes in data warehouses and datamarts which applications use to reference the various pieces of data.

[0044] According to the invention in a further regard, formulated assumptions for aspects of operation of the invention include that facts can only exist at one level of granularity, that the intersection points at any resulting row or rows on fact, and that hierarchical groupings are two-dimensional in nature. Data movement outside the data enhancement layer 110 and other portions of the supporting platform may support a push-pull relationship between the transactional and outcomes measurement layer. Extractions from source-specific to outcomes

measurement may bypass the transactional layer but may be ultimately required to feed back to support user-defined groupings. In terms of data movement of source records 118 and other data objects received or generated by the invention, functional requirements for data transport include a channel or facility for pulling data from data source 102 and set of datamarts 112, and push data to the transactional data store 130 and other repositories.

[0045] In terms of schema for ancillary physical structures according to embodiments of the invention, the translation matrix 124 may define or process at least the following functions or combinations:

# TABLE 6

Dimension and fact combinations.

Dimension and fact to source.

Dimension and fact to incident types.

Dimension and fact to incident with factors to events.

Groups to represent source-specific groups.

Groupings to represent warehouse-derived groupings – may include groups.

Ability to be represented as outcome measurement.

[0046] Once the source records 118 have been processed according to rules 120 or other paradigms and the enhanced data grouping 122 has been generated and stored to an appropriate one or more of the set of datamarts 112, according to embodiments of the invention a systems administrator, researcher or other end user may run queries against the set of datamarts 112 via query engine 114. The end user may execute those actions for instance using a user interface 116 such as an OLAP, SQL or other query or user interface, for instance using a graphical user interface interfacing to query engine 114. As an example, the end user may run a report against one or more of the set of datamarts 112 using query engine 114 and user interface 116 to

formulate a query against hospital inpatient records to ask, for instance: How many patients admitted to the hospital last year exhibited blood glucose levels above 200, along with positive detection of A1C hemoglobins?

That query might serve to detect persons having diabetes or at risk for diabetes, whether or not they were admitted or treated for that condition. Similarly, as another example a hospital administrator or other end user might execute a query against one or more of the set of datamarts 112 to determine average patient reimbursements or billings for all cardiac or oncology patients admitted in the last month. Other queries or reports are possible. According to embodiments of the invention in another regard, the complex of the set of datamarts 112, query engine 114 and user interface 116 as well as other elements or resources may together be referred to as data warehouse 136, although implementations may vary.

[0048] Overall data processing according to an embodiment of the invention is illustrated in Fig. 7. In step 702, processing may begin. In step 704, patient incident or other data may be collected at a data source 102, such as a hospital, laboratory or other site or facility. In step 706, the resulting source records 118 may be transmitted to the staging database 104 or other intermediate destination. In step 708, the source records 118 may be preprocessed, formatted or otherwise treated to permit or enhance downstream communication or processing. In step 710, the source records 118 may be conditioned by conditioning engine 106, for instance for storage in OLAP or other storage platforms.

[0049] In step 712, the data enhancement layer 110 may apply rules 120 to source records 118 or aggregations of source records 118 and other information. In step 714, data enhancement layer 110 may generate an enhanced data grouping 122. In step 716, the enhanced data grouping 122 may be stored to transactional data store 130 or elsewhere. In step 718, the

enhanced data grouping 122 and other data may be imported to the set of datamarts 112. In step 720, a systems administrator, analyst or other end user may run a report off of one or more of the set of datamarts 112, for instance to analyze disease, drug efficacy, therapeutic, demographic or other trends. In step 722, the results of any report or query may be viewed and re-queried if desired. In step 724, processing may repeat, return to a prior point or end.

[0050] The foregoing description of the invention is illustrative, and modifications in configuration and implementation will occur to persons skilled in the art. For instance, while the invention has generally been described in terms of a single data enhancement layer 110, in embodiments multiple enhancement layers may be employed. Similarly while the invention has generally been illustrated in terms of one data source 102 communicating data to the data enhancement layer 110 and other system stages, in embodiments multiple data sources may communicate a variety of source records and other information to the data enhancement layer 110 and other components.

[0051] Similarly, while the invention has in embodiments been described as processing and enhancing medical or clinical data, in embodiments data of other types may be received and treated. The scope of the invention is accordingly intended to be limited only by the following claims.